

I hereby certify that this paper is being deposited with the United States Postal Service as Express Mail in an envelope addressed to: Commissioner for Patents, Washington, D.C. 20231, on

Date: 8-14-01 *James H. Frank*  
Express Mail Label No.: EL70988462 US

## SEARCHING TOOL AND PROCESS FOR UNIFIED SEARCH USING CATEGORIES AND KEYWORDS

This application claims the benefit of priority of EPO application Serial No. 00402311.5, filed August 18, 2000, and entitled Searching Tool And Process For Unified Search Using Categories And Keywords.

### Field of the Invention

The invention relates to the field of information retrieval, and more specifically to displaying results to a search query, as well as navigating in databases and inputting requests to databases. It particularly applies to searches on the Internet.

### BACKGROUND OF THE INVENTION

Throughout the present specification, the word "site" or "internet site" refers to a number of documents connected by links, with a given entry point. A directory is the result of indexing a number of sites or documents and of classifying these into categories; categories are therefore subsets of the directory, which are usually defined in a manual operation. Such categories are often organized in a tree to facilitate navigation among categories; one may also use categories organized in a directed acyclic graph, that is a graph with a plurality of paths to the same category. A search engine is a tool for searching among documents, usually embodying automatic indexing of the documents.

A number of searching tools exist for searching and retrieving information on the Internet. Alta Vista Company proposed an Internet search site with a request box where the user may input keywords for retrieving information. The language of the search may be restricted. A box is provided that allows the user to select related searches; the related searches actually display phrases or sequences of words, which contain the current request as a substring. For instance, if the request inputted by the user reads: /greenhouse effect/ (in the rest of this specification, the request will be marked by //), related searches could offer the following choices:



consists in selecting a number of the most important terms from the documents comprised in the cluster, and in presenting them to the user. A second preferred method is to replace the important terms with important phrases, where a phrase is as a sequence of one or more words. This document provides a solution to the problem of dynamically clustering documents retrieved from a database by a search engine.

US-A-5 463 773 discloses the building of a document classification tree by recursive optimization of keyword selection function. There is provided retrieval means for extracting keywords when a document data is inputted, and outputting a classification for the document data, the classification being selected among the classification decision tree. For extracting keywords, this document suggests extracting keywords defined by word sequences. A learning process is suggested for building automatically a document classification tree on the basis of the extracted keywords.

US-A-5 924 090 proposes searching among documents, and mapping the keywords of the documents among static categories. Categories are therefore predefined in a manual process. The use of categories makes it possible to access documents included in the categories that are mapped to the categories. In this document, a search engine provides the results of a query, the results are mapped onto the static categories, and relevant categories are displayed to the user as search folders. When a search folder is selected by the user, the documents included in the search folder, that is, the documents mapped onto the corresponding category, are displayed to the user. A series of search folders is displayed any time a search is carried out, the search folders being those static categories into which a number of documents retrieved were mapped.

US-A-5 963 965 discloses a method where relevant sets of phrases are automatically extracted from text-based documents in order to build an index for these documents. These phrases are then grouped together in clusters to form a plurality of maps which graphically describe hierarchical relationships between the clusters, and can be used to extract relevant portions of the documents in answer to the user selecting one of these clusters.

US-A-5 991 756 describes a method according to which search queries may be applied to a set of documents organized in a hierarchy of categories, and where the

user is presented in response with a subset of these categories which contain the documents relevant to the query.

WO-A-98 49637 suggests organizing results of a search into a set of most relevant categories. In response to a search, the search result list is processed to dynamically create a set of search result categories. Each of the search result categories is associated with a subset of the records within the search result list having common characteristics. Categories are then displayed as folders.

The prior art information retrieval methods and processes have a number of shortcomings. Fixed or static categories actually provide a representation of the world – a set of documents – at a given time point and for a given field of the art. They may need updating, or adapting to new types of documents, when and if the set of documents is completed by new documents, especially by documents in a new field of the art. While static categories may therefore represent accurately the expertise of the human being who defined them, they are in fact limited to this expertise. In addition, any set of categories is limited by the amount of human work needed for completing categories and mapping entries of the database to the categories.

Clusters formed of keywords may provide a dynamic vision of the word. However, they do not provide an easily browsable tool, and do not allow the user to navigate easily and freely among documents.

Category searches are adapted to searching among sites. Keyword searches are more adapted to searching among separate textual documents. Therefore, there is a need for an information retrieving process and tool that enables a user to navigate not only among fixed categories, but also among keywords.

#### SUMMARY OF THE INVENTION

The present invention thus proposes a searching tool and process enabling its user to freely navigate among categories and keywords, in a friendly and transparent fashion. The invention combines the advantages of a set of human-made categories, notably expertise in a given field, together with the advantages of a keyword search, notably the ability to process and handle documents outside of said given field. The present invention provides a tool that is well adapted to searching among a database of sites and separate documents or pages.

More specifically, the present invention provides a process for searching a database of entries, including the steps of: a) providing a database of entries, at least part of said entries being mapped to a set of categories, at least part of said entries being associated with keywords, b) in response to a query of a user, selecting categories among said set of categories according to the entries returned by said query, c) dynamically selecting keywords associated to the entries returned by said query, and d) displaying to the user said selected categories and said selected keywords.

In a particular embodiment of the process, the categories are organized in a tree or directed acyclic graph structure. A keyword may preferably be a sequence of words or a sequence of stemmed words.

Selected categories and keywords may be displayed similarly, or separately. In a particular embodiment of the process, a new query is started when a user activates one of said displayed categories and keywords. This step of activating may include refining the said query of the user to the said category or keyword. This step of activating may also include excluding from the said query of the user a displayed category or keyword.

In another embodiment of the process, a list of entries returned by the query is displayed to the user. One may then display in said list a category to which at least an entry of said list is mapped. When the user selects this category in the list, the entries included in the selected category may be displayed. In this case the entries included in said category may be ranked before they are displayed.

Additionally, when categories are hierarchically organized, the step of displaying may include displaying categories of different hierarchical levels. The invention also provides a searching tool, including a search server for receiving queries from users and transmitting results to users, a database of entries, at least part of said entries being mapped to a set of categories, at least part of said entries being associated with keywords; wherein the search server includes means for searching the database and for selecting categories among said set of categories according to the entries returned by said query, means for dynamically selecting keywords associated to the entries returned by said query, and wherein the results transmitted to the users comprise said selected categories and said selected keywords.



dynamically obtained from the documents provided in response to the query. The invention allows the user to refine the search using predefined categories; in addition, displaying keywords allows the user to navigate more easily among the results of the query, without being limited to the fixed categories.

In the rest of the specification, the invention is disclosed in reference to its preferred embodiment; the database covers the World Wide Web, and includes Internet sites as well as Web documents. The tool according to the invention allows the database to be searched thanks to a HTTP server.

More generally, the invention may apply to any database where entries are at least partially mapped to predetermined categories, and may be associated with keywords. Mapping is usually a manual operation, although it is possible to use any automatic process. Textual entries may easily be associated with keywords, e. g. by automatically indexing the entries and selecting keywords. In this case, automatic mapping to categories may be carried out based on keywords.

Figure 1 is a display of a searching tool according to the invention. A request box 1 is displayed to the user, for inputting a number of keywords for a search or query. In the example of figure 1, the inputted search is again /greenhouse effect/. In a way known per se, the search may be limited to part of the database, in the example of figure 1 due to line 3. In the example of figure 1, the search is not limited, and the "World Wide Web" selection appears in bold on line 3. The "OK" button 5 permits the user to start the search or query. The current search path is displayed to the user under the request box. The use of the search path will be explained in reference to figures 2 and 3.

In response to the search, a number of documents or sites are returned. Retrieval of documents, that is, selection of sites or documents among a database of indexed or partially indexed documents or sites, may be carried out in any way known in the art. It is notably possible to use an inverted index, such as the AltaVista Search Developer's Kit, sold by AltaVista Company. More specifically, a query inputted in the query box by the user is parsed into an internal representation, which is then translated into a request applied to the inverted index. This request is formulated according to the features supported by the inverted index. Usually supported features include ranking, boolean searches, phrase searches, stemming, proximity searches, etc.







the resulting documents are a subset of the category. When a query is refined by a keyword, the resulting documents need not be associated statically with the keyword.

The difference between keywords and categories appears in Figures 1-5. Along the search, categories of the higher level disappear, as the user navigates downwards in the hierarchy of categories. This is the case where categories are organized hierarchically, e.g., in a tree or directed acyclic graph.

In case the invention is applied to a database of textual entries, there is provided an inverted index for retrieving entries. Categories are then necessarily entries of the inverted index, while keywords are not necessarily entries of the inverted index.

For instance, assume keywords are sequences of words. The keyword "fossil fuels" could then be associated to every document that contains the exact sequence of words "fossil fuels," at the time the database is built. When the keyword "fossil fuels" is selected by the user as a refinement strategy, the query may return documents not only containing "fossil fuels," but also documents containing separate occurrences of the words "fossil" and "fuels." Examples of algorithms for processing keywords – e.g. thanks to stems or synonyms – are given below.

For dynamically deriving the keywords from the entries returned by the query, one may use any process known in the art. In this respect, the keywords may include words or a sequence of words. As displayed in figure 1, it is preferred that the keywords be comprised of sequences of words. One may for instance, in a first phase of operation conducted prior to the queries, derive from all documents in the database a set of relevant sequences of words, e.g., using the algorithm described by Y. Choueka in "Looking for Needles in a Haystack or Locating Interesting Collocational Expressions in Large Textual Databases" (Conference on User-Oriented Content-Based Text and Image Handling, MIT, Cambridge, MA, pp. 609-623, 1988). This step forms a database of keywords. Then, in a second phase of operation conducted in response to the user's query, one may dynamically extract from this keyword database the keywords associated with the set of documents selected by the query and select the most frequently appearing ones as the set of keywords to be displayed to the user.

For enhancing retrieval effectiveness of the search engine, keywords may also include stems (or prefixes), instead of words, or sequences of stems of prefixes. For





suggests only performing stemming at some point of the refinement process. One may chose to do so after a given number of refinements.

Note that the same argument applies to refinement by categories. Refinement by a category may indeed return more documents than the original query, inasmuch as this query was extended by stemming. However, the refinement by category does provide a subset of the stemmed query.

In this example, it appears clearly that the documents returned after the query is refined are not necessarily associated to the keyword "fossil fuels." The exact way keywords are handled depends on the inverted index used to retrieve documents.

It is also possible to allow the user to select more than one category or one keyword as a refinement strategy. Selecting several refinement strategies at the same time may allow the user to restrict the number of documents more easily and faster. Figure 4 is a display obtained by refining the search to the keyword "CO2 emissions" in the display of figure 2. The number of results returned is low, six in the example. In this case, it is possible, as exemplified on figure 4, to display the list of documents or sites, without any further categories or keywords. This applies notably where the number of results is less than 10, under the assumption that the user may easily browse all answers, and need not restrict the number of hits.

In figures 3, 4 and 5, the search path provides links to previous keywords or searches. For instance, in figure 4 or 5, the user may select "Issues" in the search path, and get back to the display of figure 2.

Figure 5 is another view of the display of the searching tool of figure 1, after the search is limited to a category of the list of results. Specifically, the display of figure 5 is obtained when the user selects the category "Climate Change" in the sixth result of the list of results. As explained in reference to figure 1, in response to a user query, documents matching the query may be returned to the user, together with a description and, when available, the categories in which this document is classified. When the user selects on one of these categories in the list, the search engine initiates a new search and displays all documents contained in the category. Preferably, the documents are ranked or sorted according to the request box. In the example of figure 5, category "Climate Change" includes 122 documents, all of which are displayed in the list of results. The search path shows the path to the category. Sub-categories, that is, categories referenced in the "Climate Change," are displayed.









Figure 7 further shows a crawler 56 used for referencing web servers 58. The crawler searches for new information available on the Internet, and updates the database. In the examples given above, categories have a single attribute. In other words, categories are formed of a single tree. The current category may be embodied by a pointer in a directory of categories or by a pointer to one category in the graph of categories. The invention is not limited to this form of categories. Categories could be formed of several attributes. For instance, assume the invention is applied to a professional directory. In response to a request for restaurants, proposed categories could indicate the type of restaurant, the range of prices, the geographical area, and the like. These are independent attributes. A category may then be defined as a ordered set of attributes; each attribute is selected within a directory of possible values for this attribute. The current category would then be embodied by a set of pointers, each pointer pointing out to a specific value of an attribute in the relevant directory. The use of such categories makes it possible to refine a search based on several criteria. The search may be refined independently in each attribute of a category.

The invention was disclosed in the present description in reference to Internet searches, the results of the search being documents and web sites of the World Wide Web. The invention applies more generally to searches among any type of indexed or non-indexed database, provided a number of keywords may be associated to entries of the database. In addition, the entries of the database may be at least partially mapped into categories, for returning categories and allowing the user to refine the search. In this respect, the World Wide Web is a paradigm of a database, while indexed documents or web sites are paradigms of database entries. In the embodiment of the invention disclosed in figures 1 to 5, there is suggested to display a list of entries returned by the queries. The invention may actually be carried out without displaying this list, but simply by displaying refinement strategies to the user.

Last, the invention is not limited to the description made above. Other ways of populating databases may be used. Specific embodiments of a search tool and method according to the present invention have been described for the purpose of illustrating the manner in which the invention may be made and used. It should be understood that implementation of other variations and modifications of the

invention and its various aspects will be apparent to those skilled in the art, and that the invention is not limited by the specific embodiments described. It is therefore contemplated to cover by the present invention any and all modifications, variations, or equivalents that fall within the true spirit and scope of the basic underlying principles disclosed and claimed herein.

09929463, 0814401  
10-T80-29462660